

Storypoint Problem Exploration

September 2, 2024

1 Storypoint Prediction: Problem Exploration

1.1 Problem Statement

In modern agile development settings, software is developed through repeated cycles (iterative) and in smaller parts at a time (incremental), allowing for adaptation to changing requirements at any point during a project's life. A project has a number of iterations (e.g. sprints in Scrum). Each iteration requires the completion of a number of user stories, which are a common way for agile teams to express user requirements.

There is thus a need to focus on estimating the effort of completing a single user story at a time rather than the entire project. In fact, it has now become a common practice for agile teams to go through each user story and estimate its "size". Story points are commonly used as a unit of measure for specifying the overall size of a user story.

1.2 Problem Formulation

Input: A string of length N that contains a story's name and description $C = \{c_1, c_2, c_3, \dots, c_n\}$. For each story, a set of text embeddings that contains features $E = \{e_1, e_2, e_3, \dots, e_m\}$ extracted from C has been provided.

Output: A natural number P associated with the story point of that user story

1.3 Dataset Information

Text Embeddings: Text embeddings are a way to convert words or phrases from text into a list of numbers, where each number captures a part of the text's meaning. The dataset has been preprocessed and converted into two kinds of text embeddings. You can choose to work with either of them or both: - **Doc2Vec:** Input strings are transformed into fixed-length vectors of size 128. These vectors capture the semantic meaning of words and their relationships within a document. - **Look-upTable:** Input strings are transformed into fixed-length vectors of size 2264. These vectors are obtained via transforming each word in the input strings into an identifier number, then padded to the length of the longest sample.

Dataset Structure & Format: Storypoint Estimation Dataset is stored in 3 folders labeled *raw data*, *look-up*, and *doc2vec*. Within each folder are 3 CSV files for training, testing, validation. Each csv file has the following columns: - **issuekey** : The unique identifier for a story. - **storypoint**: The correct number of storypoint. - An embedding column (**embedding** or **doc2vec**) contains text embedding vectors. The raw data csv will not have this and instead contain two columns with **story name** and **description**.

1.4 Exploration

1.4.1 Raw data exploration

```
[ ]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

from sklearn.feature_extraction.text import CountVectorizer
```

Output exploration

```
[ ]: # Import raw data from the CSV file

project_name = 'aptanastudio'

all_data = pd.concat([pd.read_csv('data/' + project_name + '/' + project_name +
    ↪ '_train.csv'),
                      pd.read_csv('data/' + project_name + '/' + project_name +
    ↪ '_valid.csv'),
                      pd.read_csv('data/' + project_name + '/' + project_name +
    ↪ '_test.csv')])

print('Check the shape of the dataset', all_data.shape)
```

Check the shape of the dataset (771, 4)

```
[ ]: all_data.drop(['issuekey'], axis=1, inplace=True)
all_data.head()
```

```
[ ]:                                     title \
0  add copy url actions rightclick context menu r...
1  aptana tries open new instance opening files v...
2  content assist popup shown outside screen boun...
3                php autocompletion parent methods
4  dragging image html editor create image tag

                                     description  storypoint
0  able connect appcelerator bucket drag drog cop...         5
1  divpassigned href tender issue reported tender...         8
2  see image javascript content assist popup disp...         8
3  overwrite parent method often need call new me...        13
4  dragging image editor create img tag view beha...         8
```

First, let take a look at the distribution of the story point:

Interpretation of Skewness Values:

- **Skewness > 0:** Right-skewed distribution.
- **Skewness < 0:** Left-skewed distribution.

- **Skewness = 0**: Symmetrical distribution (like a normal distribution).

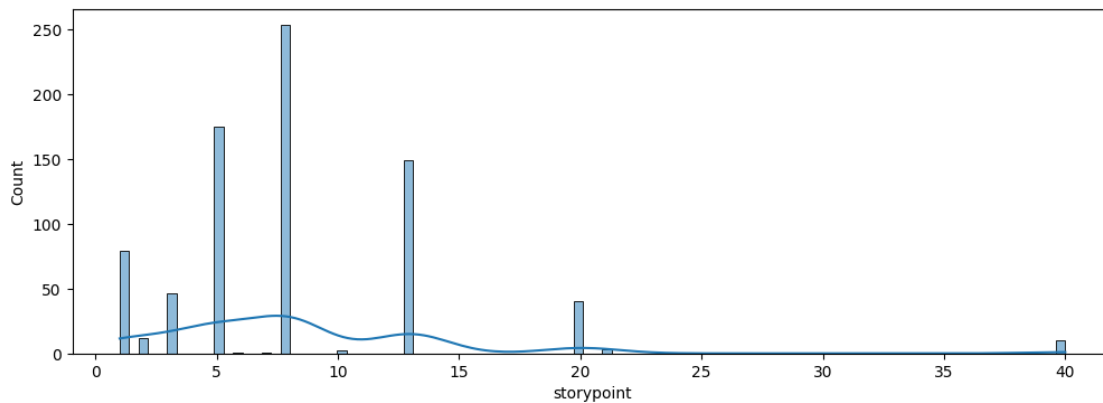
Interpretation of kurtosis: - **Leptokurtic (Kurtosis > 3)**: The distribution has heavier tails and a sharper peak than the normal distribution. Data points are more likely to produce extreme values. The distribution has a higher peak and fatter tails. - **Platykurtic (Kurtosis < 3)**: The distribution has lighter tails and a flatter peak than the normal distribution. Data are fewer extreme values compared to a normal distribution. - **Mesokurtic (Kurtosis = 3)**: The distribution has a similar kurtosis to the normal distribution, indicating a moderate level of outliers.

```
[ ]: # Draw a histogram of the story points
plt.figure(figsize=(12, 4))
plt.xticks(np.arange(0, max(all_data['storypoint']) + 1, 5))
sns.histplot(all_data['storypoint'], bins=100, kde=True)

print('Skewness:', all_data['storypoint'].skew())
print('Kurtosis:', all_data['storypoint'].kurt())
```

Skewness: 2.231210370938936

Kurtosis: 8.874730912046223



```
[ ]: tmp = pd.concat([all_data['storypoint'].value_counts(),
                    all_data['storypoint'].value_counts() / all_data.shape[0] * 100],
                    axis=1, keys=['Counts', 'Percentage (%)'])
tmp.head(20)
```

```
[ ]:
storypoint
8          253      32.814527
5          175      22.697795
13         149      19.325551
1           79      10.246433
3           46       5.966278
20          40       5.188067
```

2	12	1.556420
40	10	1.297017
21	3	0.389105
10	2	0.259403
7	1	0.129702
6	1	0.129702

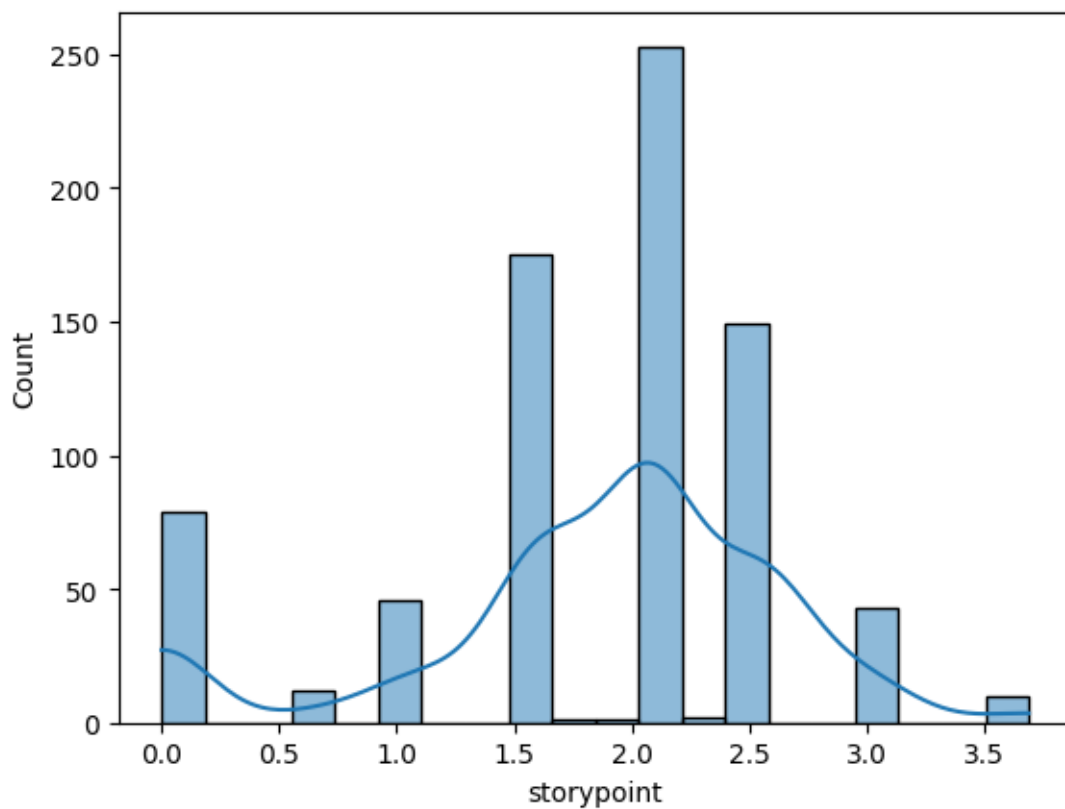
At the first sight, this data is bad. Then take a look at the statistic values, this data is even worse. Its distribution of the label is **right-skewed** and **leptokurtic**. This means if we use this to train model, the right side of the data can be the outliers and make the models become unsuable.

I will try 2 solutions: - Use log-scale on the label - Remove all the examples with label greater than a threshold (20, 30 or 40)

The first solution: logarithm magic

```
[ ]: sns.histplot(np.log(all_data['storypoint']), bins=20, kde=True)
```

```
[ ]: <Axes: xlabel='storypoint', ylabel='Count'>
```

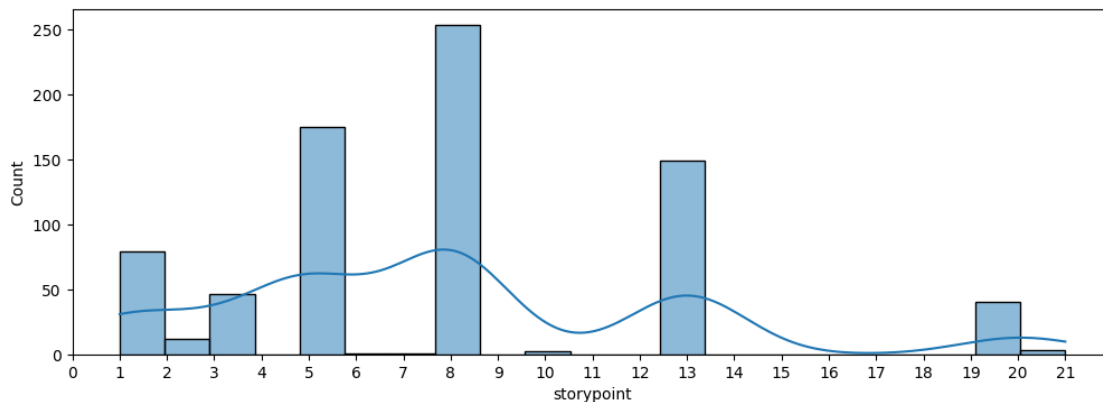


The second solution: Dismantle and Cleave

```
[ ]: threshold = 21 # This threshold means that we will take all the examples with
    ↪ story points less than or equal to 21

new_data = all_data[all_data['storypoint'] <= threshold]
plt.figure(figsize=(12, 4))
plt.xticks(np.arange(0, max(new_data['storypoint']) + 1, 1))
sns.histplot(new_data['storypoint'], bins=threshold, kde=True)
print('Fitered percentage: ', round(1 - new_data.shape[0] / all_data.shape[0],
    ↪ 2) * 100, '%')
```

Fitered percentage: 1.0 %



Input exploration The input of this problem is 2 texts: title and description. First we will find some statistics:

```
[ ]: title_lengths = all_data['title'].apply(lambda x: len(x.split(' ')))
print('Title analysis:')
print(' - Mean length:', round(title_lengths.mean()))
print(' - Min length:', title_lengths.min())
print(' - Max length:', title_lengths.max())

description_lengths = all_data['description'].apply(lambda x: len(x.split(' '))
    ↪ if type(x) != float else 0)
print('Description analysis:')
print(' - Mean length:', round(description_lengths.mean()))
print(' - Min length:', description_lengths.min())
print(' - Max length:', description_lengths.max())
```

Title analysis:

- Mean length: 6
- Min length: 1
- Max length: 18

Description analysis:

- Mean length: 50
- Min length: 3
- Max length: 1174

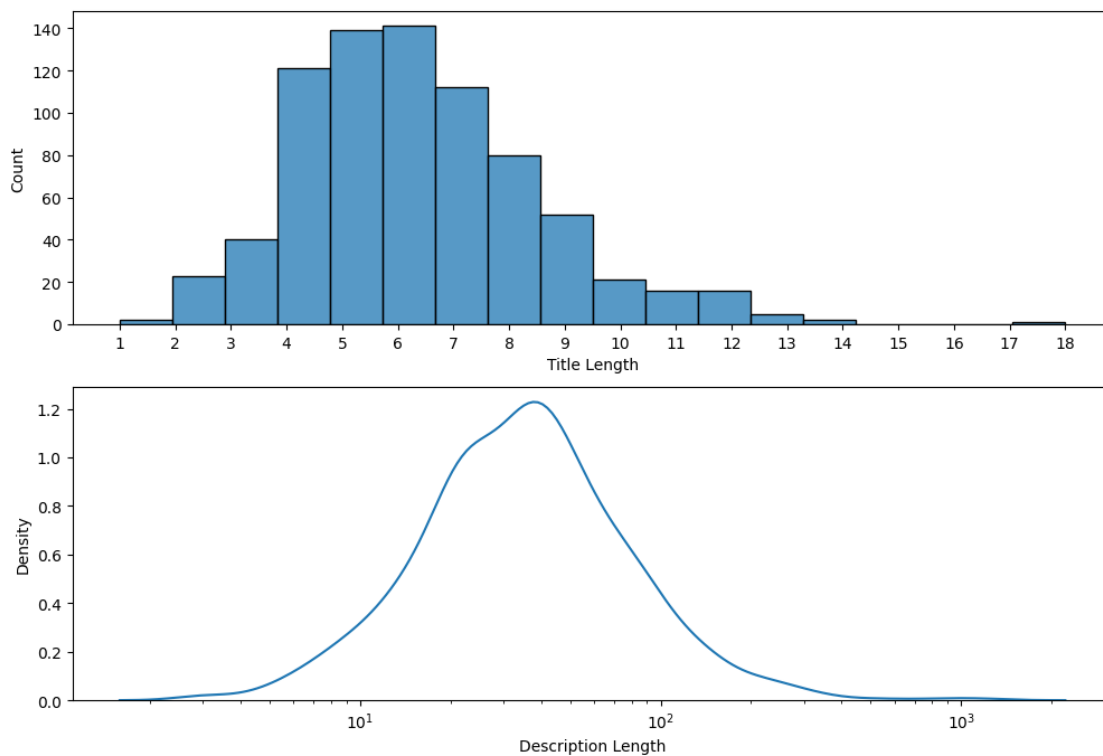
Plot the histogram of the title length and KDE of the description length (exclude 0):

```
[ ]: plt.figure(figsize=(12, 8))

plt.subplot(2, 1, 1)
plt.xticks(np.arange(0, max(title_lengths) + 1, 1))
plt.xlabel('Title Length')
sns.histplot(title_lengths, bins=max(title_lengths))

plt.subplot(2, 1, 2)
plt.xlabel('Description Length')
plt.xscale('log')
sns.kdeplot(description_lengths[description_lengths > 0])
```

```
[ ]: <Axes: xlabel='Description Length', ylabel='Density'>
```



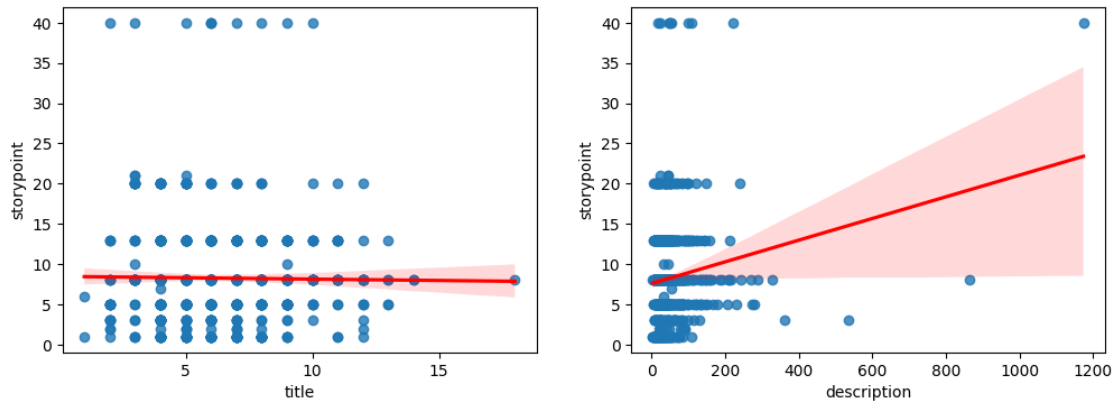
I think we should check the correlation between title length and description length:

```
[ ]: plt.figure(figsize=(12, 4))
```

```
plt.subplot(1, 2, 1)
plt.xticks(np.arange(0, max(all_data['title'].apply(lambda x: len(x.split(' ')))) + 1, 5))
sns.regplot(x=all_data['title'].apply(lambda x : len(x.split(' '))),
            y=all_data['storypoint'],
            line_kws={'color': 'red'})

plt.subplot(1, 2, 2)
sns.regplot(x=all_data['description'].apply(lambda x : len(x.split(' ')) if
type(x) != float else 0),
            y=all_data['storypoint'],
            line_kws={'color': 'red'})
```

```
[ ]: <Axes: xlabel='description', ylabel='storypoint'>
```



There is a big change in description, more description give more storypoint. However the deviation is too big, this could make the description length become a noise feature.

Let dive deeper in the input:

Title analysis:

```
[ ]: count_vectorizer = CountVectorizer()
count_vectorizer.fit(all_data['title'])

dictionary = pd.DataFrame(list(count_vectorizer.vocabulary_.items()),
                           columns=['word', 'frequency'])
dictionary.sort_values(by='frequency', ascending=False, inplace=True)
print(dictionary.shape)
dictionary.head(10)
```

```
(1589, 2)
```

```
[ ]:      word  frequency
1329      zend      1588
1362  yielding      1587
1501      yield      1586
783      yaml      1585
432   xulrunner      1584
1451   xmllang      1583
237      xml      1582
1450   xhtml      1581
1547  wrongly      1580
312      wrong      1579
```

Description analysis:

```
[ ]: count_vectorizer = CountVectorizer()
count_vectorizer.fit(all_data[all_data['description'].isnull() ==
↳False]['description'])

dictionary = pd.DataFrame(list(count_vectorizer.vocabulary_.items()),
↳columns=['word', 'frequency'])
dictionary.sort_values(by='frequency', ascending=False, inplace=True)
print(dictionary.shape)
dictionary.head(20)
```

(6726, 2)

```
[ ]:      word  frequency
2712      zzz      6725
5959   zugriff      6724
6683   zipfile      6723
3173   zipcode      6722
2037      zip      6721
4533  zfcakekohana      6720
2512      zero      6719
1494      zend      6718
5509   zadrozny      6717
3314   yvelocity      6716
4108      yui      6715
1192   youve      6714
5772   youre      6713
3218   youll      6712
415    youd      6711
5667   yields      6710
5811   yield      6709
296    yet      6708
207    yes      6707
2740   years      6706
```

Yet I don't find any thing special about the words in input except so many things are bad.

1.4.2 Solving strategies

My first intuition in this problem is that the hard part is not on the algorithm we use, it is on the **embedding** part. Therefore, in case the given embedded datasets work not properly, I will use a better embedding method which is **Bidirectional Encoder Representations from Transformers (BERT)**. Also, I will try an old way to embedding the text too: **Bag of words**.

In conclusion, I will have 4 ways to embed the text: - doc2vec (already available) - Look up (already available) - Bag Of Words - BERT

About algorithm, I will try all the regression algorithm that may give a good result:

- Ridge Regressor
- Support Vector Regressor
- Random Forest Regressor
- Gradient Boosting
- XGBoost
- Lightgbm
- Blended

Maybe, we can change the problem to the classification problem with 100 labels (desparation confirmed). In the classification problem, I will use: - Support Vector Classifier - Softmax Regression (Multinomial Logistic Regression) - Random Forest - Adaboost - XGBoost

Thanks to the libraries, the implementation of all the algorithm shrinks to its minimum form.

At last, there is still a situation that all of mentioned model don't give a good result. This gamble is thrilling (hopeless).

"But would you lose?"

Nah, I'd win.